

MENU

SEARCH

INDEX

DETAIL

BACK

2/2



JAPANESE PATENT OFFICE

PATENT ABSTRACTS OF JAPAN

(11)Publication number: 07274965

(43)Date of publication of application: 24.10.1995

(51)Int.Cl.

C12N 15/00
 // G01N 33/50
 G06F 12/00

(21)Application number: 06275336

(71)Applicant:

KOKURITSU IDENGAKU
 KENKYUSHO
 FUJITSU LTD

(22)Date of filing: 09.11.1994

(72)Inventor:

GOJIYUBORI TAKASHI
 TATENO YOSHIO
 IKEO KAZUO
 KAWANISHI YUICHI
 KAWAI MASATO

(30)Priority

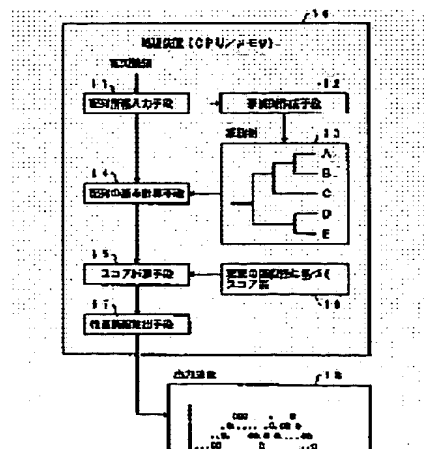
Priority number: 05283329 Priority date: 12.11.1993 Priority country: JP

(54) APPARATUS FOR EXTRACTION TREATMENT OF MOTIF OF GENE AND TREATING METHOD

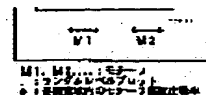
(57)Abstract:

PURPOSE: To provide an apparatus for the extraction treatment of the motif of a gene, capable of mechanically (automatically) extracting a regularity characteristic to a sequence and enabling the specification of a gene function from gene sequence information and to provide a method for the extraction treatment.

CONSTITUTION: Plural gene sequences are inputted, each sequence is weighted from alignment data on the sequences based on a phylogenetic tree, the score of each site is calculated from the weight of each sequence and the analogy of amino acids and a part having large score is automatically



extracted as a motif having a regularity characteristic to the sequence.



LEGAL STATUS

[Date of request for examination] 08.08.1997

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision of rejection]

[Date of extinction of right]

Copyright (C); 1998 Japanese Patent Office

[MENU](#)[SEARCH](#)[INDEX](#)[DETAIL](#)[BACK](#)

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平7-274965

(43) 公開日 平成7年(1995)10月24日

(51) IntCl. ⁹	識別記号	庁内整理番号	F I	技術表示箇所
C 1 2 N 15/00				
// G 0 1 N 33/50	P			
G 0 6 F 12/00	5 1 0 B	7608-5B	C 1 2 N 15/ 00	Z
		9281-4B		

審査請求 未請求 請求項の数 9 O L (全 19 頁)

(21) 出願番号 特願平6-275336

(22) 出願日 平成6年(1994)11月9日

(31) 優先権主張番号 特願平5-283329

(32) 優先日 平5(1993)11月12日

(33) 優先権主張国 日本 (J P)

(71) 出願人 593206872

国立遺伝学研究所長

静岡県三島市谷田1111番地

(71) 出願人 000005223

富士通株式会社

神奈川県川崎市中原区上小田中1015番地

(72) 発明者 五條堀 孝

静岡県三島市谷田1111番地 国立遺伝学研究所内

(72) 発明者 館野 義男

静岡県三島市谷田1111番地 国立遺伝学研究所内

(74) 代理人 弁理士 伊東 忠彦

最終頁に続く

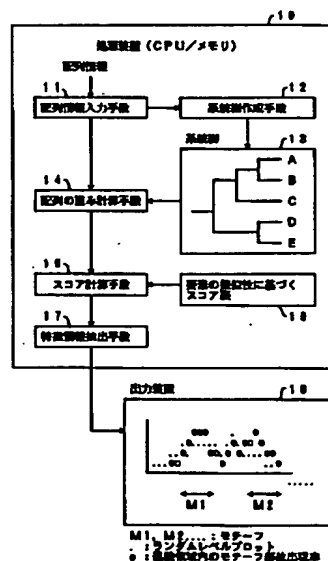
(54) 【発明の名称】 遺伝子のモチーフ抽出処理装置及び処理方法

(57) 【要約】

【目的】 遺伝子配列情報から遺伝子機能を特定する、配列に特徴的な規則性を抽出する遺伝子のモチーフ抽出処理装置および処理方法に関し、遺伝子配列情報をもとに機械的に（自動的に）モチーフを抽出することを目的とする。

【構成】 入力した複数の遺伝子配列のアライメントデータから、進化系統樹に基づく各配列への重み付けを行い、各配列の重みとアミノ酸の類似性から各部位のスコアを計算し、スコアが大きい部分を配列に特徴的な規則性であるモチーフとして自動抽出する。

本発明の構成例を示す図



【特許請求の範囲】

【請求項1】 遺伝子配列情報から遺伝子機能を特定する特徴的な規則性を抽出する処理装置であって、複数の遺伝子配列に関する進化系統樹（13）の枝の長さから各配列の重みを計算する配列の重み計算手段（14）と、

配列の各部位毎に、その部位において出現する配列要素の類似性の度合を示すスコアを前記重みを用いて計算するスコア計算手段（15）と、

計算されたスコアに基づいて遺伝子配列における特徴的な規則性を有する部分をモチーフとして抽出する特徴情報抽出手段（17）とを備えた、
遺伝子のモチーフ抽出処理装置。

【請求項2】 請求項1記載の遺伝子のモチーフ抽出処理装置において、

複数の遺伝子配列のアライメントデータをもとに遺伝子配列間の相違度に基づく進化系統樹（13）を作成する系統樹作成手段（12）を更に備えた、
遺伝子のモチーフ抽出処理装置。

【請求項3】 請求項1又は2記載の遺伝子のモチーフ抽出処理装置において、

前記スコア計算手段（15）は、前記配列の重み計算手段（14）によって計算された各配列の重みと、予め配列要素の種類に応じて求められている配列要素間の類似性情報（16）とに基づいてスコアを計算する手段を有する、
遺伝子のモチーフ抽出処理装置。

【請求項4】 請求項1～3のうちいずれか一項記載の遺伝子のモチーフ抽出処理装置において、

前記特徴情報抽出手段（17）は、前記スコア計算手段（15）によって計算されたスコアの値が所定の閾値または設定された閾値を超えた場合に、その部位をモチーフ部位として抽出し出力する手段を有する、
遺伝子のモチーフ抽出処理装置。

【請求項5】 請求項4記載の遺伝子のモチーフ抽出処理装置において、

前記特徴情報抽出手段（17）は、所定の連続領域幅又は設定した連続領域幅でモチーフ部位の出現率を計算し、その値が所定のランダムレベル又は設定したランダムレベルを超えた場合に、その連続領域をモチーフ領域として抽出すると共に、隣合う領域がともにモチーフ領域の場合は、それらの領域を1つのモチーフ領域として出力する手段を更に有する、
遺伝子のモチーフ抽出処理装置。

【請求項6】 請求項5記載の遺伝子のモチーフ抽出処理装置において、

少なくともモチーフ部位の出現率やランダムレベルの値をプロットしてグラフ表示することにより前記特徴情報抽出手段（17）からの特徴情報を出力する出力手段（18）を更に有する、

遺伝子のモチーフ抽出処理装置。

【請求項7】 計算機により遺伝子配列情報から遺伝子機能を特定する特徴的な規則性を抽出する処理方法であって、

抽出対象となる複数の遺伝子配列のアライメントデータを入力する処理過程と、

それをもとに進化系統樹を作成する処理過程と、

その系統樹における枝の長さから各配列の重みを計算する処理過程と、

計算された各配列の重みと、予め配列要素の種類に応じて求められている配列要素間の類似性情報とに基づいて、各部位のスコアを計算する処理過程と、

計算されたスコアの値が所定の閾値または設定された閾値を超えた場合に、その部位をモチーフ部位として抽出する処理過程とを有する、
遺伝子のモチーフ抽出処理方法。

【請求項8】 請求項7記載の遺伝子のモチーフ抽出処理方法において、

所定の連続領域幅又は設定した連続領域幅でモチーフ部位の出現率を計算し、その値が所定のランダムレベル又は設定したランダムレベルを超えた場合に、その連続領域をモチーフ領域として抽出すると共に、隣合う領域がともにモチーフ領域の場合は、それらの領域を1つのモチーフ領域として出力する処理過程を更に有する、
遺伝子のモチーフ抽出処理方法。

【請求項9】 請求項8記載の遺伝子のモチーフ抽出方法において、

少なくともモチーフ部位の出現率やランダムレベルの値をプロットしてグラフ表示する前記領域をモチーフとして抽出する処理過程からの特徴情報を出力する処理過程を更に有する、
遺伝子のモチーフ抽出処理方法。

【発明の詳細な説明】

【0001】

【産業上の利用分野】 本発明は遺伝子のモチーフ抽出処理装置及び処理方法に係り、特に与えられた複数の遺伝子配列情報の比較からそれらの配列間の保存部位であるモチーフを抽出する遺伝子のモチーフ抽出処理装置及び処理方法に関する。近年の遺伝子工学の進歩に伴い、DNA配列やアミノ酸配列で表現される遺伝子配列情報データベースが急増している。また、ヒトゲノム計画などのように、特定の生物の遺伝子配列を全て解明しようという試みが世界的規模で行われており、遺伝子配列情報は今後も急激に増加することが予想される。

【0002】 これらの遺伝子配列の中には、配列情報は明らかになっているが、その機能や構造に関しては未知であるものも多い。このような遺伝子の機能や構造を、その配列情報から予測するために有効な方法として、配列に特徴的な規則性であるモチーフの検索がある。そのために、配列が既知のものから多くのモチーフを抽出す

る技術が必要とされる。

【0003】

【従来の技術】従来、遺伝子配列において遺伝子機能を特定する、配列に特徴的な規則性を示すモチーフは、実験や文献での報告に基づいて決定されてきた。このようなモチーフを登録したデータベースとして、PROSITEが知られている。ところで、一般に、遺伝子配列の中で、機能的に重要な部位（サイト）は変わりにくいことが知られている。このことを利用すれば、複数の遺伝子配列の比較から、保存領域としてモチーフを抽出することができる。しかしながら、従来、遺伝子配列の比較からモチーフを抽出する手法は確立されていない。

【0004】

【発明が解決しようとする課題】実験等により人為的にモチーフを決定するのは、大変な作業である。そこで、遺伝子配列の比較からモチーフを機械的に抽出することができれば、遺伝子機能の解明等に有効な多くの情報を得ることができると考えられる。しかし、単に複数の遺伝子配列の各部位を比較し、各部位の類似性を調べていく手法を採った場合、次のような問題がある。

【0005】つまり、抽出対象とする複数の遺伝子配列情報が特定の種類の生物に偏った場合、抽出しようとする規則性に偏りが生じる。例えば、人間の遺伝子配列情報、猿の遺伝子配列情報、馬の遺伝子配列情報、・・・等の高等生物の遺伝子配列情報が多数あり、それより下等な生物の遺伝子配列情報が少ない配列情報群について、各部位の類似性からモチーフを抽出しようとした場合、類似性の高い部分が進化においてあまり変化していない保存領域であるとは必ずしも認定することはできず、モチーフとして抽出する保存領域の認定が誤りが生じる可能性がある。この逆の場合も同様である。

【0006】本発明は上記問題点の解決を図り、複数の遺伝子配列情報をもとに、機械的に（自動的に）モチーフを抽出することを目的とする。

【0007】

【課題を解決するための手段】図1は本発明の構成例を示す図である。図1において、10はCPU及びメモリ等からなる処理装置である。配列情報入力手段11は、モチーフの抽出対象となる複数の遺伝子配列のアライメントデータを入力する手段である。系統樹作成手段12は、配列情報入力手段11によって入力した複数の遺伝子配列のアライメントデータをもとに遺伝子配列間の相違度に基づく進化系統樹13を作成する手段である。なお、系統樹13は、例えば古生物学的な情報等を用いて予め作成しておくようにしてもよい。

【0008】配列の重み計算手段14は、系統樹13の枝の長さから各配列の重みを計算する手段である。スコア計算手段15は、配列の各部位毎に、その部位において出現する配列要素の類似性の度合を示すスコアを、配列の重みと、予め配列要素の種類に応じて求められてい

る要素の類似性に基づくスコア表16に基づいて計算する手段である。

【0009】特徴情報抽出手段17は、計算されたスコアに基づいて遺伝子配列における特徴的な規則性を有する部分をモチーフとして抽出し、ディスプレイやプリンタ等の出力装置18に出力させる手段である。特に、特徴情報抽出手段17は、スコア計算手段15によって計算されたスコアの値が所定の閾値または設定された閾値を超えた場合に、その部位をモチーフ部位として抽出する。又、特徴情報抽出手段17は、所定の連続領域幅又は設定した連続領域幅でモチーフ部位の出現率を計算し、その値が所定のランダムレベル又は設定したランダムレベルを超えた場合には、その連続領域をモチーフ領域とし、隣合うモチーフ領域を1つのモチーフ領域とする。これらの計算結果等が出力装置18へ出力される。出力装置18は、モチーフの出現率やランダムレベルの値等をプロットしてグラフ表示を出力する。

【0010】

【作用】本発明の遺伝子のモチーフ抽出処理装置は、マルチプルアライメントデータを入力データとし、各部位においてアライメントデータを構成する配列中で高度に保存されているアミノ酸をモチーフとして出力する。ただし、進化的に近縁な配列が存在することによる、アミノ酸の出現頻度の偏りを補正するために、アライメントデータに基づき系統樹13を作成し、系統樹13の枝長や形から、各遺伝子配列に対する重み付けを行う。更に、性質の似たアミノ酸の出現を許容するために、アミノ酸の類似性に基づいて計算されたスコア表16を用いて、各部位でのスコアを計算する。ここで求められたスコアが高いほど、その部位では、アミノ酸が高度に保存されていることを示す。

【0011】更に、モチーフ部位を抽出するために、スコアの閾値を設定する操作を行う。又、ここで設定した閾値を超えるスコアを示した部位をモチーフ部位として抽出する。そして、モチーフ領域を限定するために、領域幅とランダムレベルとを設定する操作を行う。ここで設定した領域幅内でのモチーフ部位の出現率がランダムレベルを超える値を示した場合、その領域をモチーフ領域とみなす。又、隣合うモチーフ領域は、1つのモチーフ領域とみなす。

【0012】

【実施例】以下、図面を参照しつつ、本発明の実施例をアミノ酸配列で表される遺伝子配列情報を例にして説明する。図2は本発明の実施例の処理フローチャートである。以下の説明における処理(a)～(k)は、図2に示す処理(a)～(k)に対応する。

【0013】図3に示す5本の遺伝子配列A～Eからなるアライメントデータを考える。アルファベット一文字がひとつのアミノ酸に対応し、配列長の*はギャップを表す。

10

20

30

40

50

(a) アライメントデータ入力

配列情報入力手段11は、図3に示す配列A~Eのアライメントデータを入力する。配列情報がよく似た配列が多数存在する場合、各部位を代表するアミノ酸をその出現頻度から求めると偏りが生ずる。そこで、以下の処理では、入力したアライメントデータから系統樹13を作成し、系統樹13の枝長や形をもとに、各遺伝子配列に対する重み付けの計算を行う。その計算結果を用い、各配列に対して重み付けを行うことで、偏りを補正する。

【0014】 (b) 系統樹作成

系統樹作成手段12による系統樹13の作成には、例えばUPG (Unweighted Pair-Group Clustering) 法を用いる。他の作成方法を用いてもよい。本実施例では、具体的には系統樹の作成を以下のように行う。まず、アライメントデータをもとに、遺伝子配列間の相違度を求める。相違度は2本ずつの配列を組にして、それら配列間のアミノ酸の置換数として計算される。計算式はアミノ酸置換数を求める時に一般的に使われる次の式(1)を用いる。

【0015】

$$K = -\log(1-p) \quad \dots \text{式(1)}$$

ここで、Kはアミノ酸置換数、pは2本の配列間で異なるアミノ酸を持つ部位の割合である。また、ギャップを含む部位については計算から除外する。式(1)により、全ての2本の配列の組、即ち(A, B), (A, C), ..., (A, E), (B, C), ..., (C, D), (C, E), (D, E)の組について相違度を計算する。また、この相違度を V_{AB} , V_{AC} , ..., V_{DE} と表すと、相違度 V_{AB} , V_{AC} , ..., V_{DE} の中で最小のものを選び、その組を結び付ける。この例では、配列Dと配列Eが結び付けられる。この相違度を枝の長さとする。

【0016】次に配列Dと配列Eを一つのグループとし、これらと他の各配列との相違度を同様に式(1)により計算する。例えば、配列D, Eと配列Aとの相違度 $V(DE)A$ は、 $V(DE)A = (V_{AD} + V_{AE}) / 2$ で求められる。同様に、 $V(DE)B$, $V(DE)C$ についても計算し、これらと、前に求めた V_{AB} , V_{AC} , ... から V_{DE} を除いたものの中から、最小の値を持つものを選ぶ。この例では、配列Aと配列Bの相違度 V_{AB} が最小であり、これらが2番目にグループ化される。以下、同様にグループ化と相違度の計算を行い、その結果をもとに系統樹13を作成する。

【0017】図4(A)は、図3に示すアライメントデータをもとに構築された系統樹13の例であり、図中の括弧内の数字は系統樹13中の各枝の長さを表している。

(c) 各遺伝子配列に対する重み付けの計算

次に、配列の重み計算手段14は、作成された系統樹13の各枝の長さをもとに、各枝に重みを付与する。各枝

に与える重みは、その枝から分岐した配列の本数で枝の長さを割ることにより求める。

【0018】図4(A)の系統樹13を例に説明する。枝1は長さ0.158であり、枝1からは配列A、配列B、配列Cの3本の配列が分岐している。従って、枝1の重みは $0.158 / 3 = 0.053$ と求められる。同様に、全ての枝の重みを求めると、次のようになる。

枝1の重み = $0.158 / 3 = 0.053$

10 枝2の重み = $0.903 / 2 = 0.452$

枝3の重み = $0.367 / 2 = 0.184$

枝4の重み = 0.745

枝5の重み = 枝6の重み = $0.378 / 2 = 0.189$

枝7の重み = 枝8の重み = 0.000

こうして求めた各枝の重みをもとに各配列の重みを計算する。

【0019】各配列に付与する重みは、系統樹13の根(root)から遡った時に通る枝の重み合計として求める。図4(A)の系統樹13の例では、次のようになる。配列Aは、系統樹13上で枝1、枝3、枝5を通る。各枝に与えられた重みは、それぞれ0.053, 0.184, 0.378である。従って、配列Aの重みは、これらの合計で0.615と求められる。同様に、全ての配列に対する重みを計算する。更に、全ての配列の重みの合計を求め、その合計で各配列の重みを割り、重みの合計が1になるように標準化する。図4

(B)は、図4(A)の系統樹13から求めた各配列A~Eの重みを示す。

【0020】 (d) 各アミノ酸ごとの重み計算

次に、配列の重み計算手段14は、各配列A~Eの標準化された重みをもとに、各部位におけるそれぞれのアミノ酸の重みを求める。そして、部位毎に出現するアミノ酸の重みを、そのアミノ酸が現れるすべての配列の重みの合計として求める。

【0021】図3に示すアライメントデータと図4

(B)の配列の重みをもとに説明する。第1番目の部位では、配列AはQ(グルタミン)、配列BはL(ロイシン)、配列CはE(グルタミン酸)、配列Dと配列EはS(セリン)のアミノ酸がそれぞれ出現している。従って、第1番目の部位では、アミノ酸Qの重みには配列Aの重み0.210が与えられ、同様にアミノ酸Lには0.210、アミノ酸Eには0.272、アミノ酸Sには配列Dと配列Eの重みの和で0.308の重みがそれぞれ与えられる。その他のアミノ酸は、第1番目の部位では重み0となる。

【0022】同様に、第2番目の部位では、配列A, B, Cにアミノ酸V(バリン)が出現し、配列D, Eにアミノ酸A(アラニン)が出現している。この部位におけるアミノ酸Vの重みは、配列A, B, Cの重みの和で、

0.210+0.210+0.272=0.692 となる。第2番目の部位におけるアミノ酸Aの重みは、
0.154+0.154=0.308 となる。

【0023】同様にして、全ての部位において、各アミノ酸の重みを計算する。図5は、以上のようにして計算した第1番目の部位から第10番目の部位までのアミノ酸の重みを示す。なお、図5では、小数点以下第4桁までの計算結果を示している。図6は、配列の重み計算手段14が行う、図2に示す処理(d)の一実施例を示すフローチャートである。

【0024】図6中、ステップ31は、配列の重み計算手段14の、即ち、処理装置10のCPU内の、部位位置カウンタを初期化する。ステップ32は、配列の重み計算手段14の、即ち、処理装置10のCPU内の、遺伝子配列番号カウンタを初期化する。ステップ33は、現部位位置での現遺伝子配列のアミノ酸を、配列の重み計算手段14の、即ち、処理装置10のメモリ内の、現アミノ酸種格納領域に格納する。以下の説明で、現部位位置、現遺伝子配列、現アミノ酸種等は、夫々現在注目している部位の位置、現在注目している遺伝子配列、現在注目しているアミノ酸種等を指す。ステップ34は、現遺伝子配列の重みを、配列の重み計算手段14の、即ち、処理装置10のメモリ内の、現部位位置の現アミノ酸種のスコア格納領域に格納されているスコアに加算して格納する。ステップ35は、遺伝子配列番号カウンタに1を加える。

【0025】ステップ36は、遺伝子配列番号カウンタの値が遺伝子数より大きいかなかを判定し、判定結果がNOであれば、処理はステップ34へ戻る。他方、ステップ36の判定結果がYESであれば、ステップ37が
30 現部位位置の各アミノ酸のスコアの合計が1になるように標準化処理を行う。次に、ステップ38は、各部位でのスコアを計算し、ステップ39は、部位位置に1を加*

$$\begin{aligned}
 S_1 = & D(S, S) \times S(S) \times S(S) \\
 & + D(S, L) \times S(S) \times S(L) \\
 & + D(S, E) \times S(S) \times S(E) \\
 & + D(S, Q) \times S(S) \times S(Q) \\
 & + D(L, S) \times S(L) \times S(S) \\
 & + D(L, L) \times S(L) \times S(L) \\
 & + D(L, E) \times S(L) \times S(E) \\
 & + D(L, Q) \times S(L) \times S(Q) \\
 & + D(E, S) \times S(E) \times S(S) \\
 & + D(E, L) \times S(E) \times S(L) \\
 & + D(E, E) \times S(E) \times S(E) \\
 & + D(E, Q) \times S(E) \times S(Q) \\
 & + D(Q, S) \times S(Q) \times S(S) \\
 & + D(Q, L) \times S(Q) \times S(L) \\
 & + D(Q, E) \times S(Q) \times S(E) \\
 & + D(Q, Q) \times S(S) \times S(Q) \quad \dots \text{式(2)}
 \end{aligned}$$

ここで、 S_1 は第一番目の部位のスコア、 D (アミノ酸 50 1, アミノ酸2) はアミノ酸1とアミノ酸2のスコア表

*える。ステップ40は、部位位置カウンタの値がアライメントデータ長より大きいかなかを判定し、判定結果がYESであれば処理が終了する。他方、ステップ40の判定結果がNOであれば、処理はステップ32へ戻る。

【0026】(e) 各部位でのスコア計算

配列によっては、性質の類似したアミノ酸への置換が起こっている場合があるが、このような場合でも、機能的に保存されていることが多い。そこで、スコア計算手段15は、このような部位をモチーフとして抽出するために、アミノ酸間の物理化学的類似性に基づくスコア表をもとに各部位のスコアを計算する。

【0027】アミノ酸の類似性に基づくスコア表16aは、予め各アミノ酸の物理・化学的性質をもとに求められているものであって、各アミノ酸の組の置換頻度や性質の違いの程度を示す距離に基づいて、各アミノ酸の組に対して付与された値を持つテーブルである。例えば、グリシン(G)と他のアミノ酸との組のスコアは、次のような値が付与されている。ただし、この場合には便宜上各スコアが100倍されている。

【0028】

グリシン (G) - グリシン (G)	100.0
- アラニン (A)	74.0
- セリン (S)	75.7
- ロイシン (L)	0.0
.....	

このようなスコア表については、種々のものが知られているので、ここでの説明はこの程度にとどめる。

【0029】スコア表16aの値を加味して計算した各部位のスコアは、その値が大きいほど、その部位ではアミノ酸が「保存的」であることを示している。例えば、図5の第1番目の部位のスコアは、次の式(2)で求められる。

16aから得た類似度、S（アミノ酸）はその部位におけるアミノ酸の重み（図5）である。

【0030】図7は、スコア計算手段15が行う、図2に示す処理(e)の一実施例を示すフローチャートである。図7中、ステップ51は、スコア計算手段15の、即ち、処理装置10のメモリ内の、現部位位置のスコア格納領域を初期化する。ステップ52は、スコア計算手段15の、即ち、処理装置10のメモリ内の、アミノ酸種カウンタを初期化する。ステップ53は、スコア計算手段15の、即ち、処理装置10のメモリ内の、比較アミノ酸種カウンタを初期化する。ステップ54は、現アミノ酸種と現比較アミノ酸種間の類似度を、アミノ酸類似度スコア表16aを参照して得る。ステップ55は、 $S_i = S_i + D(A_1, A_2) \times S(A_1) \times S(A_2)$

なる計算を行う。ここで、 S_i は現部位位置(i番目)のスコア、 A_1 は現アミノ酸種、 A_2 は現比較アミノ酸種、 $D(A_1, A_2)$ は現アミノ酸種と現比較アミノ酸*

(部位01~05)	0.5183	0.7744	0.5677	0.8198	0.4881
(部位06~10)	0.9328	0.4940	0.8683	0.3165	0.3580
(部位11~15)	0.9311	0.3834	0.4072	0.3611	0.6114
(部位16~20)	0.6937	0.5976	0.5699	0.5574	0.5010
(部位21~25)	0.3880	0.6168	0.5530	0.5739	0.6296
(部位26~30)	0.7718	0.3473	0.3772	0.6956	1.0000
(部位31~35)	0.9841	0.9646	1.0000	0.9149	0.8891
(部位36~40)	1.0000	0.6916	0.7864	0.7804	0.7903
(部位41~45)	0.5830	0.6021	0.7753	0.5654	0.6976
(部位46~50)	0.9037	0.6428	0.8303	0.9542	0.7105

(g) 閾値設定

特徴情報抽出手段17は、スコアに閾値を決定し、その閾値を超えるスコアの与えられた部位をモチーフとして抽出する。そのため、閾値を、ユーザの指定またはデフォルト値として事前に定められている値により設定する。

【0033】(h) 閾値を超える部位をモチーフとして出力

スコアの閾値が t_h の場合、特徴情報抽出手段17は、次式の条件を満たす部位をモチーフの候補として抽出する。

$$S > t_h$$

図3に示すアライメントデータについて、スコアの閾値 t_h を0.90として抽出したモチーフは、以下のとおりであった。

【0034】『30 D [LI] [IM] L [LIF] [KRH] L』

ここで、「30」は図3のアライメントデータ中におけるモチーフの先頭アミノ酸の位置が30であることを意味する。また、□ は、その部位では□内の複数のアミノ酸が出現していることを示す。即ち、抽出されたモチーフ部位は、アライメントデータ中の6番目の(Fま

*種の類似度のスコア、 $S(A_1)$ は現部位位置の現アミノ酸種のスコア、 $S(A_2)$ は現部位位置の現比較アミノ酸種のスコアを夫々示す。

【0031】ステップ56は、比較アミノ酸種を次の比較アミノ酸種に変更し、ステップ57は、全てのアミノ酸種との比較が行われたか否かを判定する。ステップ57の判定結果がNOであれば、処理はステップ53へ戻る。他方、ステップ57の判定結果がYESであると、ステップ58でアミノ酸種を次のアミノ酸種に変更する。ステップ59は、全てのアミノ酸種についてスコアの計算を行ったか否かを判定し、判定結果がYESであれば処理が終了する。他方、ステップ59の判定結果がNOであれば、処理はステップ52へ戻る。

【0032】(f) 計算結果出力

図3に示すアライメントデータについて、スコア計算手段15でスコアを計算した結果は以下のとおりであった。

たはY)、11番目(FまたはY)、30番目(D)、31番目(LまたはI)、32番目(IまたはM)、33番目(L)、34番目(LまたはIまたはF)、36番目(L)、46番目(IまたはV)及び49番目(LまたはIまたはM)の部位である。

【0035】図8は、本発明の実施例によって抽出されたモチーフ部位を示す図である。ところで、従来は、遺伝子配列に特徴的な配列パターンであるモチーフを部位としてマニュアル操作である程度までは抽出することができたが、モチーフを領域として同定することは困難であった。しかし、機能領域の同定や、祖先遺伝子の推定を行う場合、モチーフを領域として同定することは非常に重要である。そこで、本発明において、部位として抽出されたモチーフ配列を領域として同定する方法について、より詳細に説明する。

【0036】(i), (j), (k) 領域幅及びランダムレベルの設定、モチーフ部位の出現率の計算、モチーフ領域の出力

図9は、特徴情報抽出手段17が行う、図2に示す処理(i), (j), (k)の一実施例を示すフローチャートである。本実施例は、大略3つの処理からなる。第1の処理では、任意の領域幅を設定し、その領域幅内のモ

チーフ部位の出現率を求める。第2の処理では、設定した領域幅内でのモチーフ部位の出現率が十分に高いか否かを判断するためのランダムレベルを求め、ランダムレベルを越える出現率でモチーフ部位が存在する場合はその領域幅内のモチーフ部位を1つのモチーフ領域として同定する。第3の処理では、同定されたモチーフ領域が連続する場合にはそれらをまとめて1つのモチーフ領域とする。

【0037】つまり、より具体的には以下の処理S1～S6が繰り返される。

S1：モチーフ部位の抽出を行う。

S2：初期領域幅、拡張幅、最大拡張幅を設定する。また、モチーフ部位の出現率のランダムレベルを求めるための領域幅を最大拡張幅に設定する。ただし、最大拡張幅がアライメントデータ長の半分を越える場合には、ランダムレベルの領域幅をアライメントデータ長の半分の長さを越えない値に設定する。

S3：初期領域幅及びランダムレベルの領域幅の夫々でのモチーフ部位の出現率を計算してプロットする。

S4：初期領域幅でのモチーフ部位の出現率がランダムレベルの領域幅でのモチーフ部位の出現率を越えている場合には、初期領域幅を「モチーフ領域」とみなす。

S5：隣合う初期領域幅のモチーフ部位の出現率がともに「モチーフ領域」である場合には、これらを結合して1つの「モチーフ領域」とみなす。

S6：処理S4及びS5をアライメントデータの全長に渡って繰り返す。

【0038】図9に基づいてモチーフ領域の同定処理を説明するに、ステップ61はモチーフ部位の出現率を求める領域幅を設定する。ステップ62は、特徴情報抽出手段17の、即ち、処理装置10のCPU内の、部位位置カウンタを初期化する。ステップ63は、現部位位置を中心として、設定した領域幅内でのモチーフ部位の出現率を、次の式から求める。

【0039】(モチーフ部位の出現率) = (領域幅内モチーフ部位数) / (領域幅)

ステップ64は、モチーフ部位の出現率をグラフにプロットし、ステップ65は、現部位位置でのランダムレベルを計算し、グラフにプロットする。ステップ66は、部位位置に1を加え、ステップ67は、部位位置カウンタの値がアライメントデータ長より大きいかなかを判定する。ステップ67の判定結果がNOであれば、処理はステップ63へ戻る。

【0040】他方、ステップ67の判定結果がYESであると、ステップ68で特徴情報抽出手段17の、即ち、処理装置10のCPU内の、モチーフ領域フラグを初期化する。ステップ69は、部位位置カウンタを初期化する。ステップ70は、部位位置のモチーフ出現率がランダムレベルより高いかなかを判定する。ステップ70の判定結果がYESであれば処理はステップ71へ進

み、NOであれば処理はステップ75へ進む。

【0041】ステップ71は、モチーフ領域フラグが立っているか(セットされているか)否かを判定し、判定結果がYESであると、ステップ72が現部位位置を中心とした領域を現モチーフ領域に加えて伸長する。他方、ステップ71の判定結果がNOであると、ステップ73はモチーフ領域フラグを立て、ステップ74は、現部位位置を中心とした領域幅の中で最初にモチーフ部位の出現する部位位置を現モチーフ領域の開始部位とする。ステップ72又は74を行った後は、処理がステップ78へ進む。

【0042】ステップ75は、モチーフ領域フラグが立っているか否かを判定し、判定結果がNOであれば、処理はステップ78へ進む。他方、ステップ75の判定結果がYESであれば、ステップ76がモチーフ領域フラグを初期化し、ステップ77が現モチーフ領域を出力する。このステップ76を行った後は、処理がステップ78へ進む。

【0043】ステップ78は、部位位置に1を加える。又、ステップ79は、部位位置カウンタの値がアライメントデータ長より大きいかなかを判定し、判定結果がYESであれば、処理は終了する。他方、ステップ79の判定結果がNOであれば、処理はステップ70へ戻る。上記の如きモチーフ領域の同定処理を行った場合の実験結果を以下に説明する。

【0044】実験では、FLAA7A-1をプローブとして、アライメントデータを対象としてモチーフ領域の同定を行った。図10～図12は、プローブ名がFLAA7A-1、homologue本数が53、初期領域幅が21、最大拡張幅が101、アライメントデータ長が97、ランダムレベルを求めるための領域幅が41、モチーフ部位抽出時の設定値が0.90の場合の実験結果を示す。図10は、設定した領域幅におけるモチーフ部位の占める割合を示しており、「○」はモチーフ領域幅の初期値でのプロットを示し、「…」はランダムレベルのプロットを示す。尚、プロットが重なった場合は、割合の高い方を優先してプロットしてある。図11は、抽出されたモチーフ部位を示す。更に、図12は、モチーフ領域の同定処理により得られたモチーフ領域を示している。図12中、「:」はモチーフの開始位置と終了位置とを示し、「[]」はそのモチーフ部位に複数のアミノ酸が出現することを示し、「-」はその部位に任意のアミノ酸又はギャップ(即ち、モチーフ部位ではない部位)が出現することを示す。

【0045】図13～図15は、同様にしてECODH FOLG-1をプローブとして用いて得られた他の実験結果を示す。図13～図15は、プローブ名がECODH FOLG-1、homologue、本数が11、初期領域幅が11、アライメントデータ長が179、ランダムレベルを求めるための領域幅が81、モチーフ部位

10

20

30

40

50

抽出時の設定値が0.90の場合の実験結果を示す。図13及び図14は、同じアライメントデータに対するプロットを分割して示しており、図15は、モチーフ領域の同定処理により得られたモチーフ領域を示す。図13及び図14は、設定した領域幅におけるモチーフ部位の占める割合をアライメントデータと対応させて示しており、「○」はモチーフ領域幅の初期値でのプロット、即ち、設定した領域幅でのモチーフ部位の出現率を示し、「...」はランダムレベルのプロットを示す。つまり、モチーフ部位の出現率がこのランダムレベルより低い場合は、このモチーフ部位をモチーフ領域とはみなさない。尚、プロットが重なった場合は、割合の高い方を優先してプロットしてある。更に、図13及び図14中、アライメントデータの左側に示されている名前は、遺伝子配列データベースDDBJに登録されている遺伝子配列のエントリー名を示す。又、Dihydrofolate reductase signature [LIF]-G-X(4)-[LIVMF]-P-Wは、モチーフデータベースPROSITEに登録されているデータである。

【0046】図15中、左側に示されている「122」等の数字は、各モチーフ領域のアライメントデータ上での開始位置を示す。又、右側に示されている「137」等の数字は、各モチーフ領域のアライメントデータ上での終了位置を示す。図16～図18は、同様にしてHUMTRX1-1をプローブとして用いて得られた他の実験結果を示す。図16～図18は、プローブ名がHUMTRX1-1、homologue、本数が15、初期領域幅が11、アライメントデータ長が110、ランダムレベルを求めるための領域幅が51、モチーフ部位抽出時の設定値が0.90の場合の実験結果を示す。図16及び図17は、同じアライメントデータに対するプロットを分割して示しており、図18は、モチーフ領域の同定処理により得られたモチーフ領域を示す。

【0047】図16及び図17は、設定した領域幅におけるモチーフ部位の占める割合をアライメントデータと対応させて示しており、「○」はモチーフ領域幅の初期値でのプロット、即ち、設定した領域幅でのモチーフ部位の出現率を示し、「...」はランダムレベルのプロットを示す。つまり、モチーフ部位の出現率がこのランダムレベルより低い場合は、このモチーフ部位をモチーフ領域とはみなさない。尚、プロットが重なった場合は、割合の高い方を優先してプロットしてある。更に、図16及び図17中、アライメントデータの左側に示されている名前は、遺伝子配列データベースDDBJに登録されている遺伝子配列のエントリー名を示す。又、Thioredoxin family active site [STA]-X-[WG]-C-[AGV]-[PH]-Cは、モチーフデータベースPROSITEに登録されているデータである。

【0048】図18中、左側に示されている「69」等の数字は、各モチーフ領域のアライメントデータ上での開始位置を示す。又、右側に示されている「105」等の数字は、各モチーフ領域のアライメントデータ上での終了位置を示す。これにより、本発明によれば、遺伝子配列情報から、機械的に（自動的に）モチーフ領域を抽出・同定することができるので、高速にモチーフ領域の抽出・同定が可能となる。従って、大量の遺伝子配列データから新規なモチーフを発見したり、モチーフデータベースを作成することが、容易にできる。この様にして得られたモチーフ情報をもとに、未知機能の遺伝子配列の機能及び構造の予測を効率良く行うことができ、本発明を遺伝子機能の発見や機能領域の同定に利用すると非常に便利である。

【0049】以上、本発明を実施例により説明したが、本発明は上記実施例に限定されるものではなく、本発明の範囲内で種々の変形及び改良が可能であることは、言うまでもない。

【0050】

【発明の効果】以上説明したように、本発明によれば、遺伝子配列情報から、機械的に（自動的に）モチーフ領域を抽出・同定することができるので、高速にモチーフ領域の抽出・同定が可能となるので、大量の遺伝子配列データから新規なモチーフを発見したり、モチーフデータベースを作成することが、容易にでき、この様にして得られたモチーフ情報をもとに、未知機能の遺伝子配列の機能及び構造の予測を効率良く行うこともできるので、本発明を遺伝子機能の発見や機能領域の同定に利用すると非常に便利であり、遺伝子工学の発展に寄与するところが大きい。

【図面の簡単な説明】

【図1】本発明の構成例を示す図である。

【図2】本発明の実施例の処理フローチャートである。

【図3】入力したアライメントデータの例を示す図である。

【図4】図3に示すアライメントデータから求めた系統樹と配列の重みの結果を示す図である。

【図5】各部位におけるアミノ酸の重みの計算結果を示す図である。

【図6】配列の重み計算手段が行う処理の一実施例を示すフローチャートである。

【図7】スコア計算手段が行う処理の一実施例を示すフローチャートである。

【図8】本発明の実施例によって抽出されたモチーフ部位を示す図である。

【図9】特徴情報抽出手段が行う処理の一実施例を示すフローチャートである。

【図10】モチーフ領域の同定処理を行った場合の実験結果を説明する図である（その1）。

【図11】モチーフ領域の同定処理を行った場合の実験

結果を説明する図である(その2)。

【図12】モチーフ領域の同定処理を行った場合の実験結果を説明する図である(その3)。

【図13】モチーフ領域の同定処理を行った場合の実験結果を説明する図である(その1)。

【図14】モチーフ領域の同定処理を行った場合の実験結果を説明する図である(その2)。

【図15】モチーフ領域の同定処理を行った場合の実験結果を説明する図である(その3)。

【図16】モチーフ領域の同定処理を行った場合の実験結果を説明する図である(その1)。

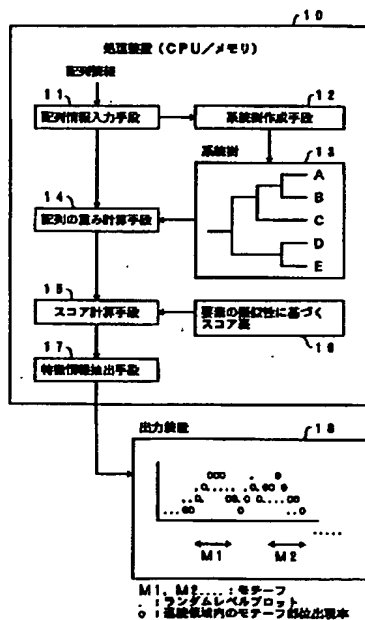
【図17】モチーフ領域の同定処理を行った場合の実験結果を説明する図である(その2)。

【図1】

【図3】

【図8】

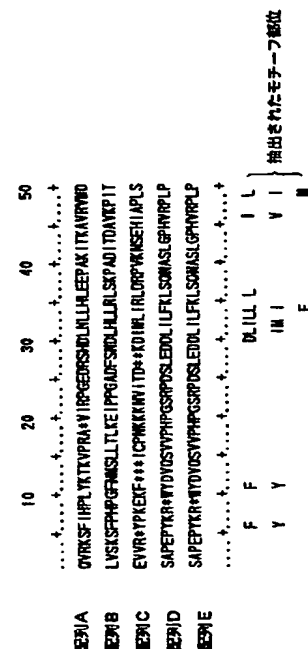
本発明の構成例を示す図



入力したアライメントデータの例を示す図

020	A	0VRKSF	1HPLTKTVPR	AV	IRGE	SD	HL	LE	EP	AK	ITKAN	VRND					
020	B	LVSKSF	PHGF	MSLL	TLKE	IPPG	AF	SN	DL	UL	SKP	ADITDA	VPIT				
020	C	EVVR	TPKE	CF	ICPK	KNN	ITD	40	IM	IRLDR	PK	SEH	IAPLS				
020	D	SAPE	TKR	WT	VS	VP	HP	GS	RP	SD	EDL	ILFL	SN	ASL	GF	HP	PLP
020	E	SAPE	TKR	WT	VS	VP	HP	GS	RP	SD	EDL	ILFL	SN	ASL	GF	HP	PLP

本発明の実施例によって抽出されたモチーフ部位を示す図



【図18】

モチーフ領域の同定処理を行った場合の実験結果を説明する図(その3)

13: [FY] :13
25: [IV]--DF-A-10GPC[KR]--[IV]-P--[LIF] :47
61: [DM]--D---[AP] :89
69: [AP]--[TIR]P[AT][LF]--[LIF]K-G---G---[LI]--[LIV] :105

【图5】

各部位におけるアミノ酸の重みの計算結果を示す図



	標準化	
ΣA	0. 015	0. 210
ΣB	0. 015	0. 210
ΣC	0. 798	0. 272
ΣD	0. 452	0. 154
ΣE	0. 482	0. 154

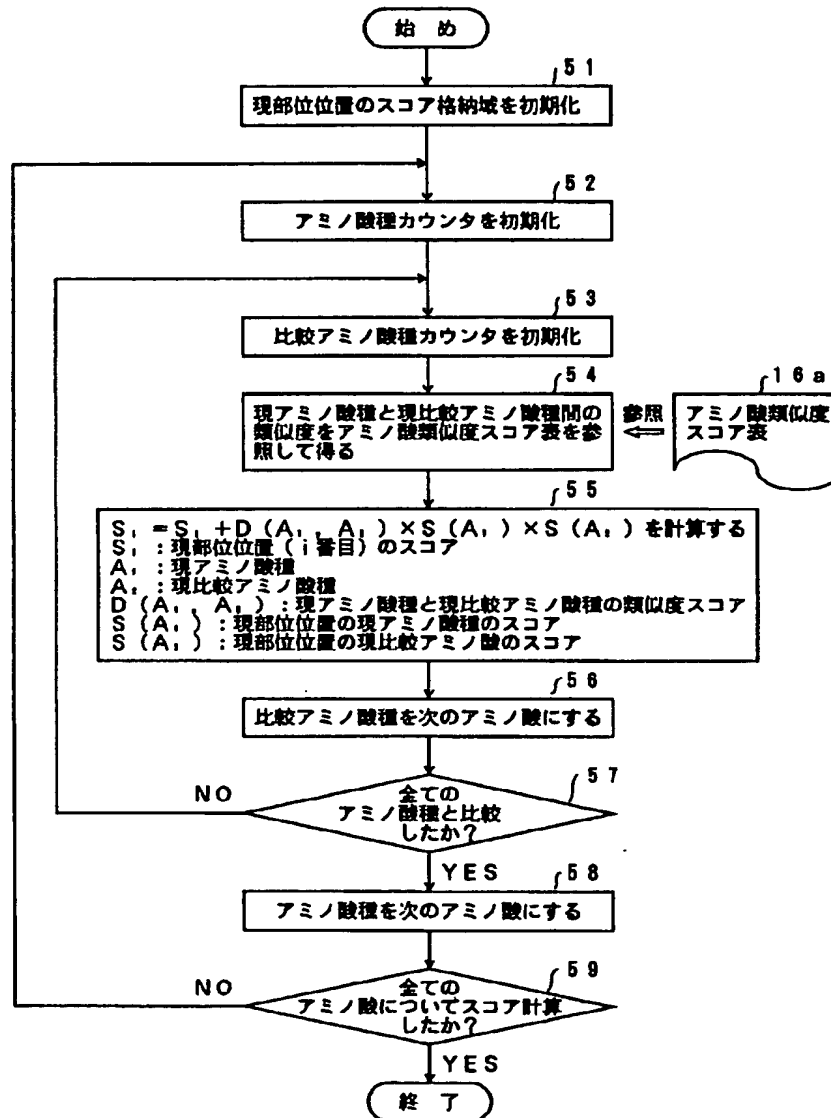
	1	2	3	4	5	6	7	8	9	0
A:	Q	V	S	K	S	F	I	H	P	J
B:	O	V	S	K	S	F	I	H	P	O
C:	J	V	S	K	S	F	I	H	P	J
D:	J	V	S	K	S	F	I	H	P	J
E:	J	V	S	K	S	F	I	H	P	J
F:	J	V	S	K	S	F	I	H	P	J
G:	J	V	S	K	S	F	I	H	P	J
H:	J	V	S	K	S	F	I	H	P	J
I:	J	V	S	K	S	F	I	H	P	J
J:	J	V	S	K	S	F	I	H	P	J
K:	J	V	S	K	S	F	I	H	P	J
L:	J	V	S	K	S	F	I	H	P	J
M:	J	V	S	K	S	F	I	H	P	J
N:	J	V	S	K	S	F	I	H	P	J
O:	J	V	S	K	S	F	I	H	P	J
P:	J	V	S	K	S	F	I	H	P	J
Q:	J	V	S	K	S	F	I	H	P	J
R:	J	V	S	K	S	F	I	H	P	J
S:	J	V	S	K	S	F	I	H	P	J
T:	J	V	S	K	S	F	I	H	P	J
U:	J	V	S	K	S	F	I	H	P	J
V:	J	V	S	K	S	F	I	H	P	J
W:	J	V	S	K	S	F	I	H	P	J
X:	J	V	S	K	S	F	I	H	P	J
Y:	J	V	S	K	S	F	I	H	P	J
Z:	J	V	S	K	S	F	I	H	P	J

【图 15】

13: A-----16---(LMDP)---(OE)---(FY)KRI : 59
 59: (LWV)GQR-T-(EN)S(NF)---LP : 68
 62: LP-----1V : 69
 07: 00---(LW)FY)---(LWV) : 122
 22: (LWV)-(LWV)---(EN)XO-(NF) : 137

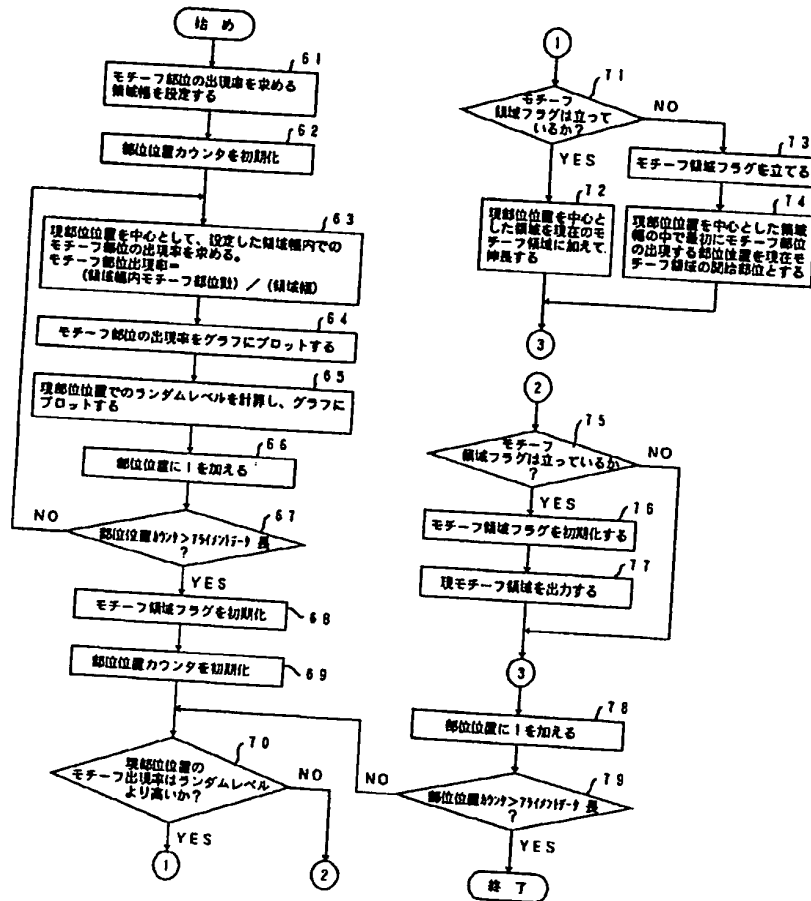
【図7】

スコア計算手段が行う処理の一実施例を示すフローチャート



【図9】

特許情報抽出手段が行う処理の一実施例を示すフローチャート



【图 17】

モチーフ領域の同定処理を行った場合の実験結果を説明する図(その2)

```

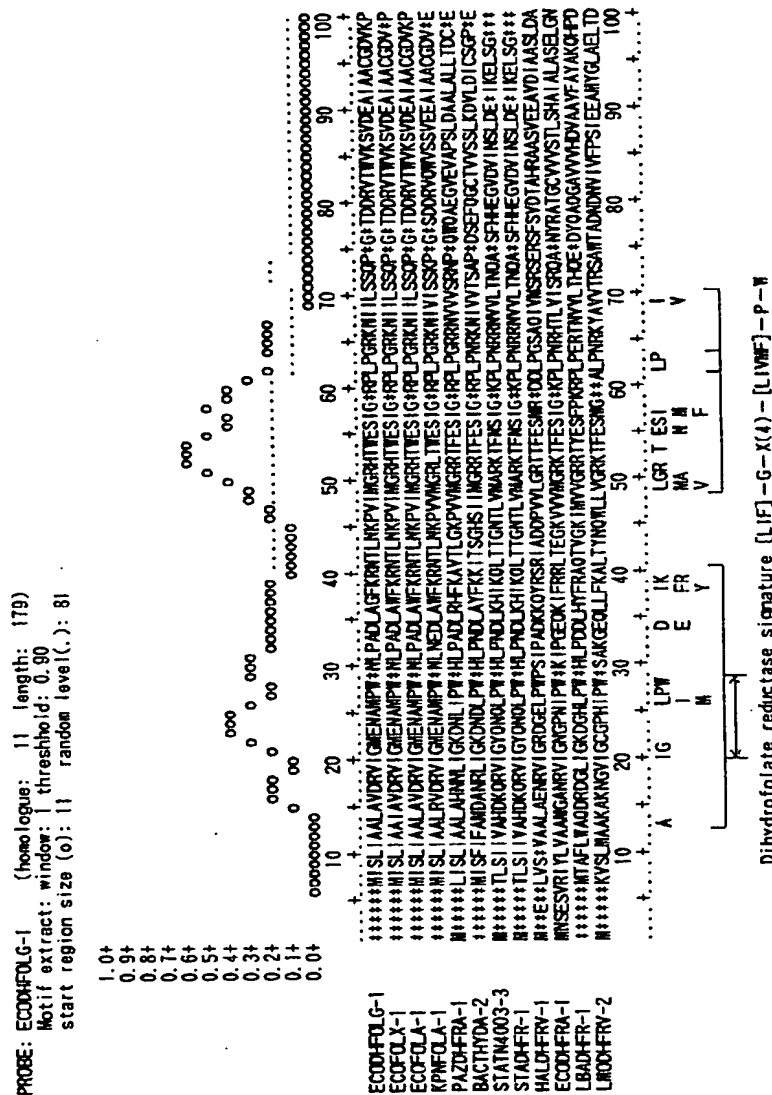
1. 0+
0. 0+
0. 0+
0. 7+
0. 6+
0. 5+
0. 4+
0. 3+
0. 2+
0. 1+oooo
0. 0+

110
.....+
HUNTXX1-1 LEAT1+HEL.V
HURTHD-1 LEAT1+HEL.V
RATTR-1 LEAT1+TEFA
CHTHD-1 LEET1+HEL.V
YSCTRIA-1 IQQ1+ASRV
YSCTRX2-1 IQQ1+ASRV
YSCTRX1-1 IQQ1+ANA
YSCTRILA-1 IQQ1+ANA
ANITRDX-1 LAMT1+HEL.V
AHATRX-1 LSOT1+ELX+
STYTRX-1 LKEFLDAN.L
ECOTTRX-1 LKEFLDAN.L
EDOTRX-1 LKEFLDAN.L
ECORIDA-1 LKEFLDAN.L
RCATRX-1 LATW1+AS.L
BACAPK1-1 LQELVNO.L
110
.....+
L L
I I
V V

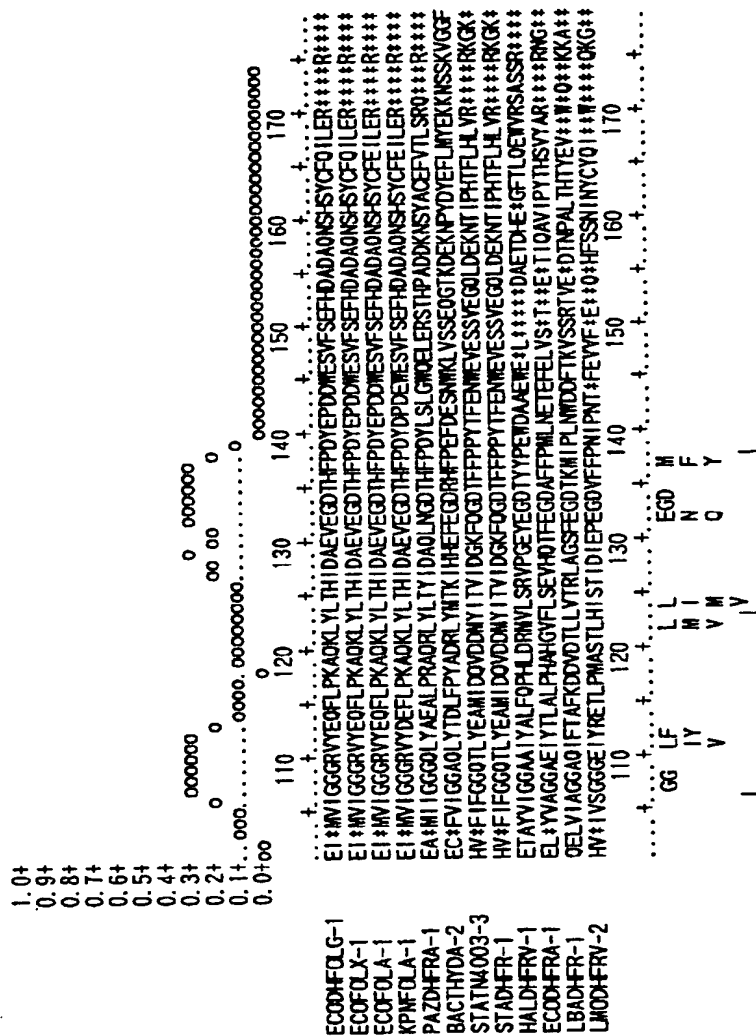
```

【図13】

モチーフ領域の同定処理を行った場合の実験結果を
説明する図（その1）

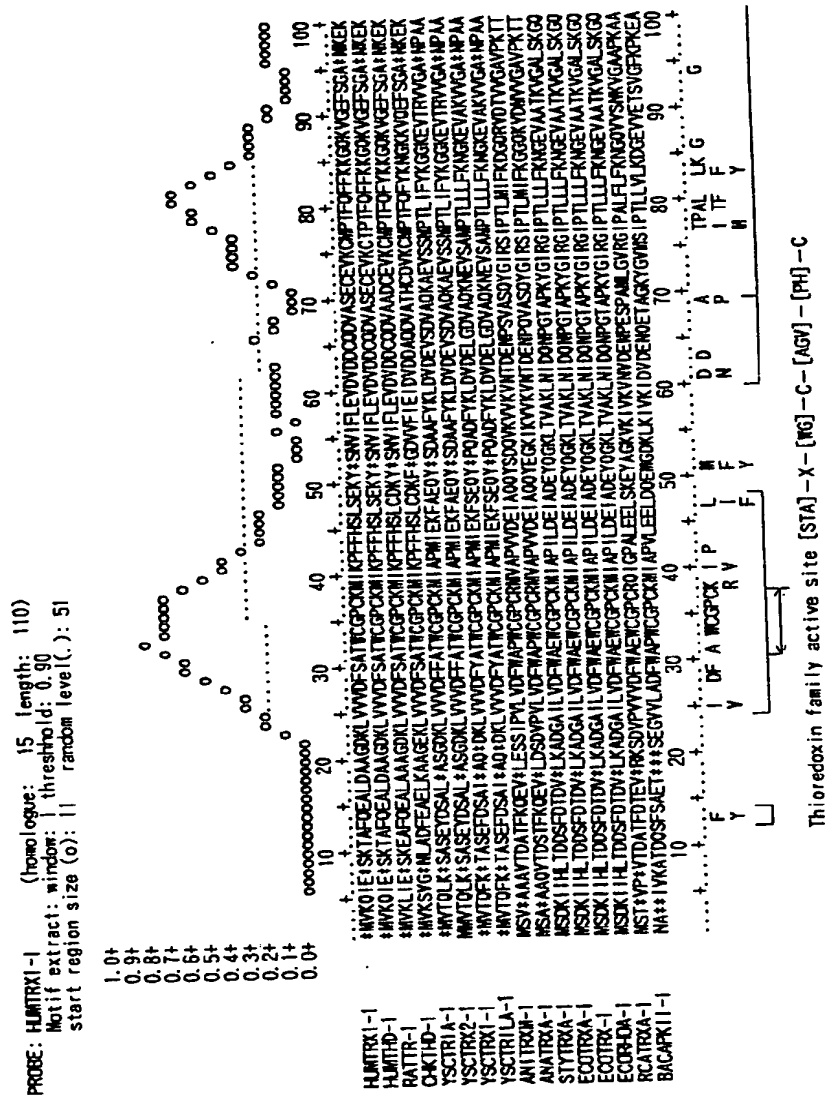


モチーフ領域の同定処理を行った場合の実験結果を説明する図(その2)



【図16】

モチーフ領域の同定処理を行った場合の実験結果を
説明する図（その1）



フロントページの続き

(72)発明者 池尾 一穂
静岡県三島市谷田1111番地 国立遺伝学研
究所内

(72)発明者 川西 祐一
神奈川県川崎市中原区上小田中1015番地
富士通株式会社内

(72)発明者 河合 正人
神奈川県川崎市中原区上小田中1015番地
富士通株式会社内